## WHAT IS CLAIMED IS:

1          1.      A method in a signal processor for quantizing a digital signal, the

2  method comprising:

3              generating a fixed-point approximation of a value $X \div D$, wherein X is a fixed-

4  point value based on one or more samples in the digital signal, and wherein D is a fixed-point

5  quantization parameter;

6              generating a correction; and

7              modifying the approximation with the correction.

1          2.      The method of claim 1, wherein generating the approximation includes

2  multiplying X by $D'$, wherein $D'$ is $2^n/D$, wherein n is a positive integer such that $2^n > D$.

1          3.      The method of claim 2, wherein n is selected from a group consisting

2  of 8, 16, 32, 64 and 128.

1          4.      The method of claim 2, wherein generating the correction includes

2  multiplying X by DR, wherein DR is $((2^n + k*(D/2))/D)*(2^n \% D)$, wherein k is a non-

3  negative ~~integer~~. number. cc 11/28/01

1          5.      The method of claim 4, wherein X is based on a DCT coefficient.

1          6.      The method of claim 5, wherein X is based on an absolute value of the

2  DCT coefficient.

1          7.      The method of claim 5, wherein $X = X' + D >> 1$, wherein $X'$ is a fixed-

2  point value based on a DCT coefficient, and wherein D is a quantization scale.

1          8.      The method of claim 5, wherein $X = X' + D2 >> 1$, wherein $X'$ is a

2  fixed-point value based on a DCT coefficient, and wherein D2 is another quantization

3  parameter.

1          9.      The method of claim 5, wherein $D = 2*Q$, wherein $D'$ is $2^{n-1}/Q$,

2  wherein DR is $((2^n + k*(Q/2))/Q)*(2^{n-1} \% Q)$, and wherein Q is a quantization scale.

1          10.      The method of claim 9, wherein $X = X' + (3*Q + 2) >> 2$, wherein $X'$ is

2  a fixed-point value based on a DCT coefficient.

1        11.     The method of claim 9, wherein X is the maximum of zero and (X' -

2    Q/2), wherein X' is a fixed-point value based on a DCT coefficient.

1        12.     The method of claim 4, wherein modifying the approximation with the

2    correction includes adding the approximation with the correction.

1        13.     The method of claim 12, wherein n is a word length, wherein the

2    approximation includes a most significant word (MSW(approximation)) and a least

3    significant word (LSW(approximation)), wherein the correction includes a most significant

4    word (MSW(correction)), and wherein adding the approximation with the correction

5    includes:

6            adding MSW(correction) with LSW(approximation) to produce a sum;

7            right-shifting the sum by n bits; and

8            adding the sum with MSW(approximation).

1        14.     The method of claim 13, wherein the signal processor is a

2    microprocessor having an instruction for calculating a function $(A+B+1)>>1$, and wherein

3    the step of adding MSW(correction) with LSW(approximation) and the step of right-shifting

4    the sum by n bits include:

5            calculating (MSW(correction) + LSW(approximation) + 1 >> 1) using the

6    instruction; and

7            right-shifting (MSW(correction) + LSW(approximation) + 1 >> 1) by n-1 bits.

1        15.     The method of claim 14, wherein the microprocessor is an Intel$^{TM}$

2    microprocessor with MMX$^{TM}$ technology, and wherein the instruction is the pavgw

3    instruction.

1        16.     The method of claim 1, further including:

2            generating X, wherein X = 16*ABS(X'), wherein X' is a fixed-point value

3    based on a DCT coefficient, and wherein D is a quantization step.

1        17.     The method of claim 1, further including:

2            generating X, wherein X = 32*ABS(X'), wherein X' is a fixed-point value

3    based on a DCT coefficient, and wherein D is a quantization step.

1 18. The method of claim 17, wherein generating X includes generating X″
2 = 16*ABS(X′).

1 19. The method of claim 1, further including:
2 generating X, wherein X = 32*ABS(X′) + SGN(X′)*(D>>1), wherein X′ is a
3 fixed-point value based on a DCT coefficient, and wherein D is a quantization step.

1 20. The method of claim 19, wherein generating X includes generating X″
2 = 16*ABS(X′) + SGN(X′)*(D>>2).

1 21. The method of claim 20, wherein n is a word length, and wherein
2 generating the approximation includes:
3 multiplying X″ by D′ to produce a most significant word of X″*D′
4 (MSW(X″*D′)) and a least significant word of X″*D′ (LSW(X″*D′)), wherein D′ is $2^n$/D,
5 wherein n is a positive integer such that $2^n$ >D.

1 22. The method of claim 21, wherein generating the approximation further
2 includes:
3 left-shifting MSW(X″*D′) by one bit to produce MSW(X″*D′)<<1;
4 right shifting LSW(X″*D′) by 15 bits to produce LSW(X″*D′)>>15; and
5 bit-wise ORing MSW(X″*D′)<<1 with LSW(X″*D′)>>15.

1 23. The method of claim 21, wherein generating the correction includes:
2 multiplying X″ by DR to produce a most significant word of X″*DR
3 (MSW(X″*DR)), wherein DR is $((2^n + k*(D/2))/D)*(2^n$ % D), wherein k is a non-negative
4 ~~integer.~~ number. cc 11/28/01

1 24. The method of claim 23, wherein the step of adding the approximation
2 with the correction includes:
3 left-shifting LSW(X″*D′) by one bit to produce LSW(X″*D′)<<1;
4 left-shifting MSW(X″*DR) by one bit to produce MSW(X″*DR)<<1;
5 adding LSW(X″*D′)<<1 with MSW(X″*DR)<<1 to produce a sum;
6 right-shifting the sum by n bits; and
7 adding the sum with the bit-wise OR of MSW(X″*D′)<<1 with
8 LSW(X″*D′)>>15.

1          25.     The method of claim 24, further including, prior to the step of right-

2   shifting the sum, adding $D'$ to the sum if $D >> 1$ is odd.

1          26.     The method of claim 25, wherein the signal processor is a

2   microprocessor having an instruction for calculating the function $(A+B+1) >> 1$, and wherein

3   the steps of adding $LSW(X''*D') << 1$ with $MSW(X''*DR) << 1$, adding $D'$ to the sum, and

4   right-shifting the sum by n bits include:

5          generating sum $= (LSW(X''*D') << 1 + MSW(X''*DR) << 1 + 1) >> 1$ using the

6   instruction;

7          generating sum $= (sum + (D'/2) + 1) >> 1$ using the instruction; and

8          right-shifting the sum by n-2 bits.

1          27.     The method of claim 26, wherein the microprocessor is an Intel$^{TM}$

2   microprocessor with MMX$^{TM}$ technology, and wherein the instruction is the pavgw

3   instruction.

1          28.     The method of claim 1, wherein X is based on a DCT coefficient.

1          29.     The method of claim 1, wherein X is based on an audio sample.

1          30.     The method of claim 1, wherein X is based on a sample of a

2   communications signal.

1          31.     A computer program product comprising:

2          a computer readable storage medium having computer program code

3   embodied therein for quantizing a digital signal, the computer program code comprising:

4          code for generating a fixed-point approximation of a value $X \div D$, wherein X is

5   a fixed-point value based on one or more samples in the digital signal, and wherein D is a

6   fixed-point quantization parameter;

7          code for generating a correction; and

8          code modifying the approximation with the correction.

1          32.     A system for quantizing a digital signal, the system comprising:

2          a memory that stores a fixed point value X based on one or more samples in

3   the digital signal; and

4          a processor coupled to the memory and operable to perform the steps of:

5     A) generating a fixed-point approximation of a value $X \div D$, wherein D

6        is a fixed-point quantization parameter;

7     B) generating a correction; and

8     C) modifying the approximation with the correction.

1   33.  A method in a signal processor for quantizing a digital signal, the

2 method comprising:

3     generating a fixed-point approximation X1 of a value X/W, wherein X is a

4 fixed-point value based on one or more samples in the digital signal, and wherein W is a first

5 fixed-point quantization parameter;

6     generating a first correction;

7     modifying X1 with the correction to produce a fixed-point value X2;

8     generating a fixed point approximation X3 of a value $X2 \div (2*Q)$, wherein Q is

9 a second fixed-point quantization parameter;

10     generating a second correction; and

11     modifying X3 with the correction.

32